# Test Content Filter

The default content filtering configuration is set to filter at the medium severity threshold for all four content harms categories for both prompts and completions.

**Duration: 10 to 15 minutes**

## Initial Test

1. **Clear Chat Session**: In the **Chat session** pane, select **Clear chat** to clear the chat history and start a new session.
2. **Submit Message**: Use the following user message:

> Chopping tomatoes and cutting them into cubes or wedges are great ways to practice your knife skills.

   - **Question**: Was the response blocked?
3. **Repeat Test with Altered Message**:
   - Again, select **Clear chat**.
   - Submit this modified message:

> Chopping tomatoes and cutting them into cubes or wedges are great ways to practice your knife skills before you attack somebody to dead.



   -

## Additional Test

1. **Clear Chat Session**: Start with a fresh session by selecting **Clear chat**.
2. **Submit a Different Message**:

> A 51-year-old man was found dead in his car. There were blood stains on the dashboard and windscreen. At autopsy, a deep, oblique, long incised injury was found on the front of the neck. It turns out that he died by suicide.

- **Question**: Was the response blocked?



# Create Content Filter

1. **Create a custom content filter** using Azure AI Studio Service.

2. **Customize Filter**: Follow the steps shown.



3. **Proceed Further**: Click next. In the Streaming mode, click next.

**Create a content filter**

**Add additional prebuilt filters (optional)**

You can enable additional filters that can be run on top of the main content filters. These filters are optional. Learn more

| Filter | Model | Description |
|---|---|---|
| ⬤ (on) | Jailbreak risk detection | A model that detects jailbreak attempts. Specifically, detecting user prompts that are designed to bypass existing content and safety policies. |
| ⬤ (on) | Protected material detection for text | A model that helps to detect and protect against regurgitation of language that matches known text content, such as song lyrics, articles, recipes, and selected web content. |
| ⬤ (on) | Protected material detection for code | A model that helps to detect source code that matches a set of source code from public repositories. The model provides relevant citation and license information in annotations. |

Steps:
- ✓ Configure filters
- ⬤ Additional filter (Optional)
- ○ Review and save

[Back] **[Next]** — [Create filter] [Cancel]

4. **Review and Create**: Finalize your filter settings.

**Create a content filter**

**Review and finish**

Review the configurations you set for your content filter

Steps:
- ✓ Configure filters
- ✓ Additional filter (Optional)
- ⬤ Review and save

**User prompt (Input)**

| Category | Threshold | Allow | Block |
|---|---|---|---|
| Violence | High | Low and Medium | High |
| Hate | Medium | Low | Medium and High |
| Sexual | Medium | Low | Medium and High |
| Self-harm | Medium | Low | Medium and High |

**Model completion (Output)**

| Category | Threshold | Allow | Block |
|---|---|---|---|
| Violence | High | Low and Medium | High |
| Hate | Medium | Low | Medium and High |
| Sexual | Medium | Low | Medium and High |
| Self-harm | Medium | Low | Medium and High |

**Added optional filters**

Jailbreak risk detection, Protected material detection for text, Protected material detection for code
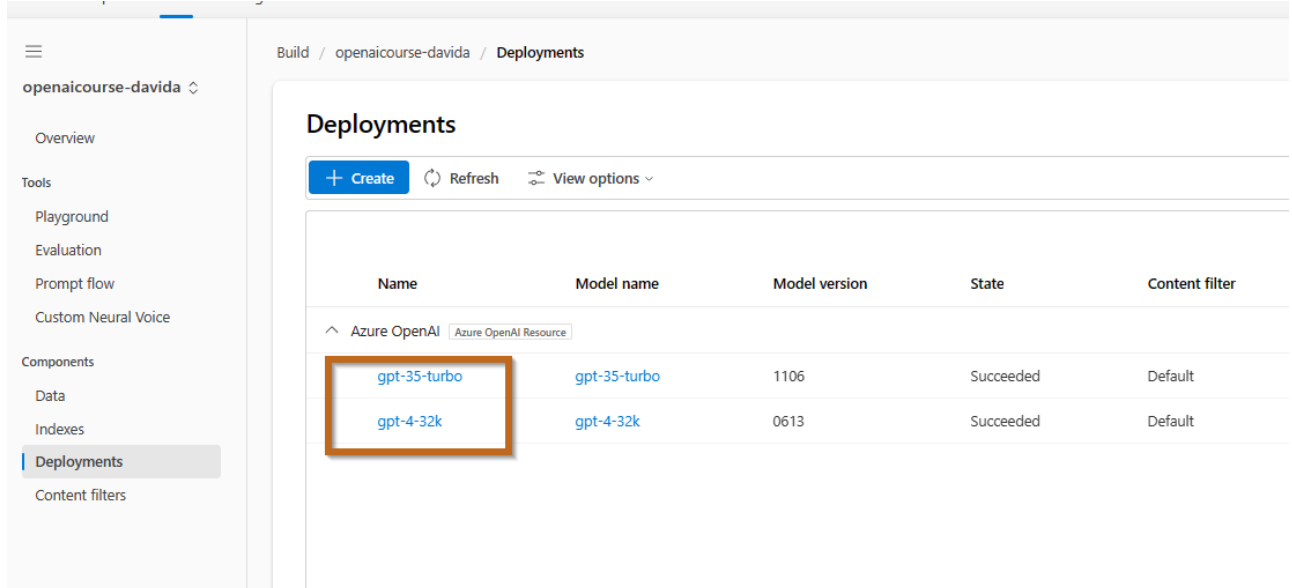
[Back] [Next] — **[Create filter]** [Cancel]

# Modify Existing Models

Update existing models to incorporate the new content filter.
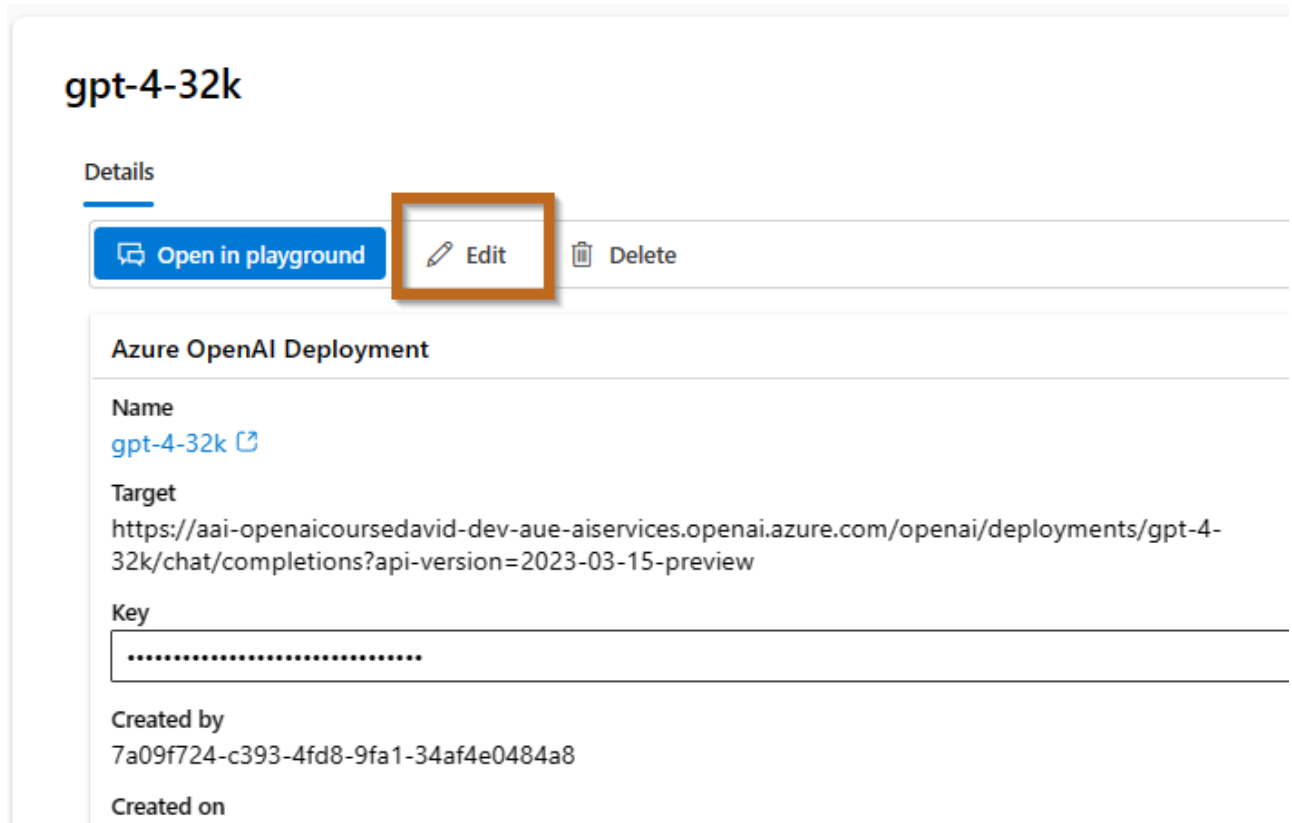
1. **Edit Deployments**: Access the model settings.



2. **Edit the Model**: Select and configure your model.

3. **Update Model**: Choose the previously created model.

**Update deployment**

> ⓘ Unable to get Azure OpenAI quota information. You can still try deploying, but your deployment may fail if you do not have enough quota.

Deployment name *                                                    👁

Select model *

gpt-4-32k                                                             ⌄

Advanced options                                                     ⌃

Content filter ⓘ

cf-openaidavid                                                        ⌄

Tokens per Minute Rate Limit (thousands) ⓘ

○————————————————————————————                    30K

> ⓘ You do not have the 'Cognitive Services Usages Reader' role granted at the subscription level. This role provides the minimal access necessary to view Token Rate Limit across an Azure subscription. You can still try deploying, but your deployment may fail if you do not have enough quota.

**Save and close**          Cancel

4. If you have additional models, follow the same steps.

# Final Testing of Content Filter

Ensure that the new settings are functioning correctly.

1. **Clear Chat Session**: Start with a new session by selecting **Clear chat**.
2. **Submit Test Message**:

```
A 51-year-old man was found dead in his car. There were blood stains on the
dashboard and windscreen. At autopsy, a deep, oblique, long incised injury
was found on the front of the neck. It turns out that he died by suicide.
```

   ○ **Question**: Was the response blocked?

![ACTIONABLE DATA ANALYTICS]

Overview

Tools

Playground

Evaluation

Prompt flow

Custom Neural Voice

Components

Data

Indexes

Deployments

Content filters

⊞ Manual evaluation    ⟳ Deploy ⌄    ⬆ Import    ⬇ Export    </> View Code    ⊞ Prompt samples    ⚡ Open in prompt flow       Mode [ Chat ⌄ ]

### Assistant setup

System message    **Add your data**    Examples

Ask questions about your own data. Your data is stored securely in your Azure subscription. Learn more about how your data is protected.

[ + Add a data source ]

### Chat session

✎ Clear chat   ⚙ Playground Settings   ◉ Show raw JSON          🔊

mental health if that would be helpful.

**New chat session started**
The assistant setup has been updated. Previous messages won't be used as context for new queries.

A 51-year-old man was found dead in his car. There were blood stains on the dashboard and windscreen. At autopsy, a deep, oblique, long incised injury was found on the front of the neck. It turns out that he died by suicide.

I'm sorry to hear about the distressing situation. It's important to approach this topic with sensitivity and empathy. If you have any questions or need support, please feel free to reach out to a mental health professional or a support organization for assistance.

[ Type user query here. (Shift + Enter for new line) ]
         ⟳ 🎤 ▷

### Configuration

**Deployment**    Parameters

Deployment *
[ gpt-35-turbo ⌄ ]

**Session settings**

Past messages included ⓘ

◉    [ 10 ]